

An Other-Race Effect for Face Recognition Algorithms

P. Jonathon Phillips, *Senior Member, IEEE*, Fang Jiang, Abhijit Narvekar, Julianne Ayyad, and Alice J. O'Toole

Abstract—Psychological research indicates that humans recognize faces of their own race more accurately than faces of other races. This “other-race effect” occurs for algorithms tested in the Face Recognition Vendor Test 2006. We report results for a Western algorithm made by fusing eight algorithms from Western countries and an East Asian algorithm made by fusing five algorithms from East Asian countries. At the low false accept rates required for most security applications, the Western algorithm recognized Caucasian faces more accurately than East Asian faces and the East Asian algorithm recognized East Asian faces more accurately than Caucasian faces. Next, using a test that spanned all false alarm rates, we compared the algorithms with humans of Caucasian and East Asian descent matching face identity in an identical stimulus set. In this case, both algorithms performed better on the Caucasian faces—the “majority” race in the database. The Caucasian face advantage, however, was far larger for the Western algorithm than for the East Asian algorithm. Humans showed the standard other-race effect for these faces, but showed more stable performance than the algorithms over changes in the race of the test faces. State-of-the-art face recognition algorithms, like humans, struggle with “other-race face” recognition.

Index Terms—face and gesture recognition, performance evaluation of algorithms and systems, human information processing

I. INTRODUCTION

THE other-race effect for face recognition has been established in numerous human memory studies [1] and in meta-analyses of these studies [2], [3], [4]. The effect for human perceivers can be summed up in the oft-heard phrase, “They all look alike to me”. This anecdote suggests that our ability to perceive the unique identity of other-race faces is limited relative to our ability to perceive the unique identity of faces of our own race. Although humans have additional social prejudices that impact our ability to recognize other-race faces [5], [6], [7], perceptual factors seem to be the primary cause of the other-race effect in humans [8], [9], [10]. These factors begin to develop early in infancy and stem from the amount and quality of experience we have with faces of different races [11]. In fact,

the other-race effect in humans can be measured in infants as a decrease in their ability to detect differences in individual other-race faces as early as three to nine months of age [11]. This occurs simultaneously with impressive gains in the ability of infants to distinguish faces of their own race. Thus, it has been argued that human deficiencies in perceiving other-race faces may be a consequence of neural feature selection processes that begin early in infant development. These processes serve to optimize the encoding of unique features for the types of faces we encounter most frequently—usually faces of our own race. The cost of this optimization is a perceptual filter that limits the quality of representations that can be formed for faces that are not well described by these features [12], [13], [14].

The rationale for testing face recognition algorithms for an other-race effect is based on the following premises. First, many face recognition algorithms include training procedures aimed at optimally representing individual faces [15], [16]. Second, the databases used for training different algorithms vary in the extent to which they represent human demographic categories. Thus, there is reason to be concerned that some of the underlying causes of the other-race effect in humans might apply to algorithms as well. Although face recognition algorithms have been tested extensively for performance stability across environmental context variables including viewpoint, illumination, and image resolution (e.g., [17], [18], [19], [20]), the question of performance stability over population demographics has received much less attention [15], [21], [22]. Furthermore, no studies have examined algorithm performance as a function of the interaction between the demographic origin of the algorithm (i.e., where it was developed) and the demographics of the population to be recognized. Understanding the stability of algorithm performance for populations of faces that vary in demographics is critical for predicting face recognition accuracy when application venues vary in their demographic structure.

In this study, performance is compared for algorithms and humans on matching identity on pairs of faces. In the identity matching task, an algorithm or human is presented with two face images and must respond with a measure of confidence to indicate whether the faces are the same person or different people. In biometrics this is referred to as a verification task. We compared the performance of an East Asian algorithm and a Western algorithm matching identity in pairs of Caucasian and East Asian faces. The East Asian algorithm was a fusion of five algorithms from East Asian countries; and the Western algorithm was a fusion of eight algorithms from Western countries. The Face Recognition Vendor Test 2006 (FRVT 2006) [23] served as the source of the algorithms that contributed to the fusions. In Experiment 1, the East Asian and Western algorithms matched face identity in all available East Asian and Caucasian face pairs from the FRVT 2006 database [23]. In this first test, we focused on a range of low false accept rates typical for security applications.

Manuscript received October 6, 2008. This work was supported by TSWG.

P. J. Phillips is with the National Institute of Standards and Technology, 100 Bureau Dr., MS 8940 Gaithersburg MD 20899, USA, (Tel. 301-975-5348), (Fax 301-975-5287), (e-mail: jonathon@nist.gov).

Fang Jiang is at the School of Behavioral and Brain Sciences, GR4.1 The University of Texas at Dallas Richardson, TX 75083-0688, USA, (Tel. 972-883-4983), (Fax 973-883-2491), (e-mail: fxj018100@utdallas.edu).

Abhijit Narvekar is at the School of Behavioral and Brain Sciences, GR4.1 The University of Texas at Dallas Richardson, TX 75083-0688, USA, (Tel. 972-883-4983), (Fax 973-883-2491), (e-mail: aln061000@utdallas.edu).

Julianne Ayyad is at the School of Behavioral and Brain Sciences, GR4.1 The University of Texas at Dallas Richardson, TX 75083-0688, USA, (Tel. 972-883-4983), (Fax 973-883-2491), (e-mail: jha011100@utdallas.edu).

A. J. O'Toole is at the School of Behavioral and Brain Sciences, GR4.1 The University of Texas at Dallas Richardson, TX 75083-0688, USA, (Tel. 972-883-2486), (Fax 973-883-2491), (e-mail: otoole@utdallas.edu).

In Experiment 2, we benchmarked the performance of the East Asian and Western fusion algorithms against the performance of humans of Caucasian and East Asian descent. This comparison was carried out using a smaller number of face pairs that allowed for a direct comparison among humans and the two algorithms. The face pairs were selected to control for demographic factors other than race. The second experiment measured performance using A' , a non-parametric statistic, which is a more general measure that characterizes performance over the full range of false accept rates. The test of humans serves as a control condition to confirm the presence of the other-race effect for the sample of faces tested. It also provides a baseline measure of human accuracy and recognition stability over a change in the race of the test population. This measure can be used to benchmark algorithm stability over demographic change.

II. EXPERIMENT 1

A. Methods

The FRVT 2006 was the source of the algorithm data and face images for this comparison [23]. The National Institute of Standards and Technology (NIST) sponsored, U.S. Government test of face recognition algorithms was open to academic and corporate researchers worldwide [23]. Algorithms in the FRVT 2006 competition were required to match facial identity in 568, 633, 560 pairs of images over five experiments on still face images. *Match pairs* consisted of two images of the same person and *non-match pairs* consisted of two images of different people. In this study, we focused on a face pairs from one of the FRVT 2006 experiments where the images varied in illumination conditions and there were a sufficient number of Caucasian and East Asian faces¹. Specifically, one image in the pair was taken under controlled illumination (e.g., under studio lighting) and the other image was taken under uncontrolled illumination (e.g., in a corridor). The uncontrolled illumination images had a resolution of 2272×1704 pixels and the controlled illumination images had a resolution of 1704×2272 pixels. Example image pairs for the East Asian and Caucasian faces appear in Figure 1.

B. Algorithms

Algorithms participating in the FRVT 2006 could be divided into those submitted by research groups from East Asia and those submitted by research groups from Western countries (Western Europe and North America). Five algorithms were submitted by research groups in East Asia (2 algorithms from China, 2 algorithms from Japan, 1 algorithm from Korea) and eight algorithms were from research groups in Western countries (2 algorithms from France, 4 algorithms from Germany, 2 algorithms from The United States). We report performance for the average of the East Asian algorithms—an *East Asian fusion algorithm*, and for the average of the Western algorithms—a *Western fusion algorithm*. (Details on the performance of individual algorithms are available elsewhere, [23]).

The task of the individual algorithms was to compare identity in pairs of face images consisting of a controlled and an uncontrolled illumination image. Identity comparisons were based on the computed similarity scores between the controlled and uncontrolled

illumination images. Specifically, the source data for Experiment 1 were based on each algorithms' matrix of similarity scores for all available East Asian face pairs age 18-35 ($n = 205, 114$; 200, 256 non-match pairs and 4, 858 match pairs) and all available Caucasian face pairs age 18-35 ($n = 3, 359, 404$; 3, 345, 592 non-match pairs and 13, 812 match pairs). These scores were extracted from the similarity score matrix computed by each algorithm for the FRVT 2006.

The algorithms from East Asian and Western countries were fused separately in a two-step process. In the first step, for each algorithm, the median and the median absolute deviation (MAD) were estimated from 6849 out of 7, 007, 032 similarity scores ($median_k$ and MAD_k are the median and MAD for algorithm k). The median and MAD were estimated from 6849 similarity scores to avoid over tuning the estimates to the data. The similarity scores were selected to evenly sample the images in the experiment. The fused similarity scores are the sum of the individual algorithm similarity scores after the median has been subtracted and then divided by the MAD. If s_k is a similarity score for algorithm k and s_f is a fusion similarity score, then $s_f = \sum_k (s_k - median_k) / MAD_k$.

C. Results

Figure 2 shows the receiver operating characteristic (ROC) curve for the performance of the algorithms. The ROC plots the trade-off between the verification rate and the false accept rate as a threshold is varied. A false accept occurs when an algorithm incorrectly states that the faces of two different people are the same person. A successful verification occurs when an algorithm correctly accepts that two faces are from the same person. The curve is plotted using a logarithmic scale on the horizontal axis to highlight performance in the range of the low false accept rates required for security applications. A classic "other-race effect" is evident. The East Asian fusion algorithm is more accurate at recognizing the East Asian faces and the Western fusion algorithm is more accurate on the Caucasian faces. There is also an advantage for East Asian faces, consistent with both uncontrolled and controlled face matching studies in the literature [21], [22], [24]. As we will see, this advantage may be primarily limited to the low false accept rate operating points common in security applications.

III. EXPERIMENT 2

In Experiment 2, we carried out a direct comparison between humans of East Asian and Caucasian descent and the East Asian and Western fusions algorithms. In the first experiment, we found an other-race effect for the East Asian and Caucasian fusion algorithms using all available pairs of East Asian and Caucasian face pairs. One limitation of using all available pairs of faces is that the people in the pairs may have differed on characteristics other than race (e.g., gender, age). In this second experiment, we used a smaller set of face pairs that were matched carefully for demographic characteristics other than race. We measured human and algorithm performance using the area under the ROC (AUC). The AUC is a general measure of performance that summarizes a ROC across all false accept rates.

A. Stimuli

Forty pairs of Asians (20 match pairs and 20 non-match pairs; 16 female and 24 male pairs) and 40 Caucasian pairs (20 match

¹The algorithm results used in this study are from the very-high resolution still face images in the uncontrolled illumination experiment (section 5.3, [23]).



Fig. 1. Example of controlled (left) and uncontrolled (right) illumination images.

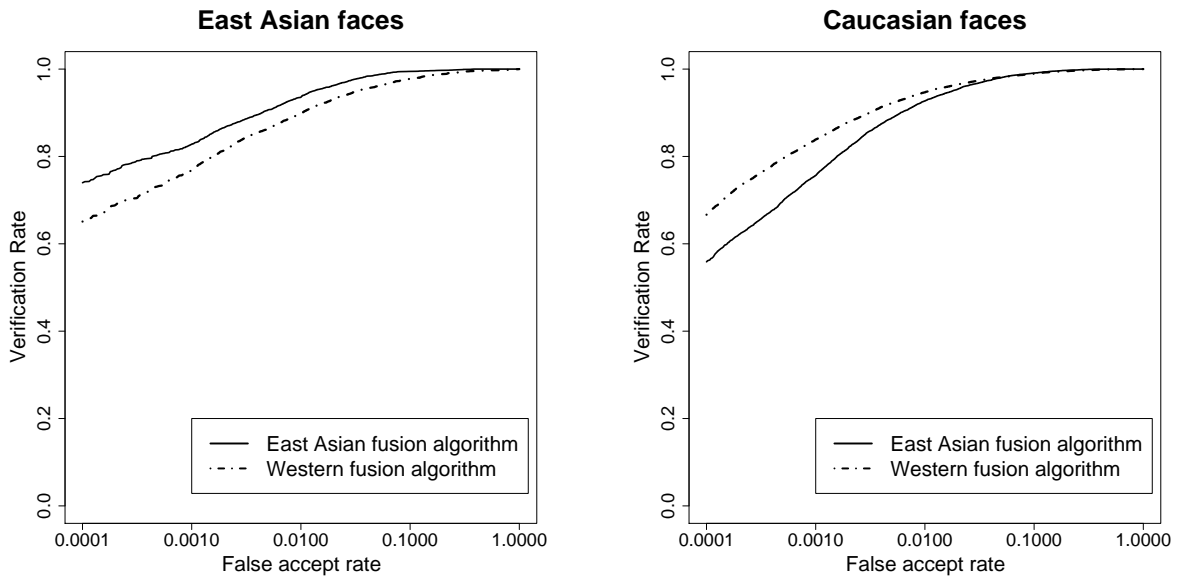


Fig. 2. ROC of the East Asian fusion and Caucasian fusion algorithms on the Experiment 1 data set. The horizontal axis is on a logarithmic scale to emphasize low false accept rates typical of security applications. The East Asian fusion algorithm is more accurate with East Asian face pairs and the Western fusion algorithm is more accurate with Caucasian face pairs. The effect is most pronounced at lower false accept rates, where security applications commonly operate.

pairs and 20 non-match pairs; 16 female and 24 male pairs) were used in the experiment. Following the procedure in the FRVT 2006 human performance studies, selected face pairs were rated as having medium difficulty; e.g., approximately half the algorithms matched the identity correctly [23], [25]. Face pairs in the experiment excluded mismatched gender and retained only young adult faces (i.e., 18-35 years old).

B. Human Experimental Methods

1) *Human Participants*: Undergraduate students from the School of Behavioral and Brain Sciences at University of Texas at Dallas volunteered to participate in these experiments in exchange for a research credit in a psychology course. A total of 26 students (19 females and 7 males) participated in the experiment. There were 16 Caucasians (11 female, 5 male) and 10 East Asians (8 female, 2 male).

2) *Procedure*: In the experiment, human participants were asked to match the identity of people in pairs of face images. On each trial, an image pair was displayed on the computer screen for 2s, followed by a prompt asking the participant to respond as follows, “1.) sure they are the same; 2.) think they are the same; 3.) do not know; 4.) think they are not the same; and 5.) sure they are not the same.” The next trial proceeded after a response was entered. Participants matched East Asian face pairs in one block of 40 trials and the Caucasian face pairs in another block of 40 trials. Half of the participants were tested with the East Asian faces first and Caucasian faces second. The remaining subjects were tested with the blocks in the reverse order.

C. Results

1) *Human Behavioral Data*: The performance of the East Asian and Caucasian participants on the East Asian and Caucasian faces was measured by tallying their responses to the match and non-match pairs. For the purpose of measuring statistical significance, we computed each subjects’ A' for discriminating match versus non-match pairs for the Caucasian and East Asian faces. The statistic A' is used commonly in the psychology literature as a non-parametric estimate of area under the ROC curve derived from human certainty data [26]. In each condition, A' was computed from the subjects’ rating responses as follows. Responses 1 and 2 were deemed “same person” judgments and responses 3, 4, and 5 were deemed “different person” judgments². For each subject, the verification rate (i.e., hit rate) was computed as the proportion of face pairs correctly judged to be “same” when the face pair were images of the same person. The false acceptance rate (i.e., false accept rate) was calculated as the proportion of face pairs that were incorrectly judged to be the same, when they were images of two different people. The statistic A' was then computed from the hit and false accept rates as

$$\frac{1}{2} + \left[\frac{(H - F)(1 + H - F)}{4H(1 - F)} \right], \quad (1)$$

where H is the hit rate and F is the false accept rate [27]³. Analogous to the AUC measure for algorithms, this formula

²Note that a reasonable alternative measure is to assign ratings 1, 2, and 3 to the category of hits. We computed all results in this alternative fashion and found the same pattern of results.

³ A' is the non-parametric version of d' . We use A' here as it approximates the area under the ROC statistic used for the algorithms. Note, however, that d' yielded the same pattern of results and statistical effects.

gives a score of 1 for perfect performance and .5 for chance performance.

The A' values for East Asian and Caucasian face identification were then used to compute a partially repeated-measures analysis of variance (ANOVA) with race of participant (East Asian or Caucasian) as a between-subjects factor and race of the face (East Asian or Caucasian) as a within-subjects factor. Evidence for the other-race effect was found in the form of a significant interaction between the race of the subject and the race of the face, $F(1, 24) = 6.15, p < .02$. The pattern of this interaction appears in Figure 3 (left side) and shows a substantial Caucasian face advantage for Caucasian participants and a slight East Asian face advantage for East Asian participants. No statistically significant effects were found for the race of the subject or for the race of the face, indicating that there was no statistical difference in the accuracy of East Asian and Caucasian participants and no statistical difference in accuracy of participants overall for the East Asian and Caucasian faces. We note however that the interaction is tilted slightly (though not significantly) in favor of accuracy on Caucasian faces. We will consider this result shortly in the Discussion in the context of the algorithm results.

D. Algorithm Methods

Next, we compared algorithm and human accuracy on these face pairs. Algorithm performance was assessed by extracting similarity scores from the East Asian and Western fusion algorithms for the same set of face pairs presented to human participants. On these face pairs, the AUC was computed for the East Asian and Western fusion algorithms on the East Asian face pairs and the Caucasian face pairs.

E. Results

The AUC values for the algorithms are shown on the right side of Figure 3 for comparison with human performance. Both the East Asian and Western fusion algorithms were more accurate with the Caucasian face pairs than with the East Asian face pairs. However, the accuracy advantage for Caucasian face pairs is far larger for the Western fusion algorithm than for the East Asian algorithm. This result is consistent with an other-race effect, but one that is superimposed on a Caucasian face advantage. Thus, the results differ from human performance in the general advantage seen for the Caucasian faces, but are similar to human performance in the interaction seen between the demographic origin of the algorithm and the race of faces in the test set. It is also worth noting that the Western and East Asian algorithms show a larger difference in performance for the two test populations than the humans. This indicates that the performance of the algorithms is less stable over race change than the performance of humans.

IV. CONCLUSION

The primary conclusion of this work is that demographic origin of face recognition algorithms and the demographic composition of a test population interact to affect the accuracy of the algorithms. At its core, this finding indicates that algorithm performance varies over changes in population demographics. As noted, the variability of algorithm performance over changes in viewing parameters such as illumination and pose has been well studied previously. The results of the present study indicate that stability over demographics should also be included in measures

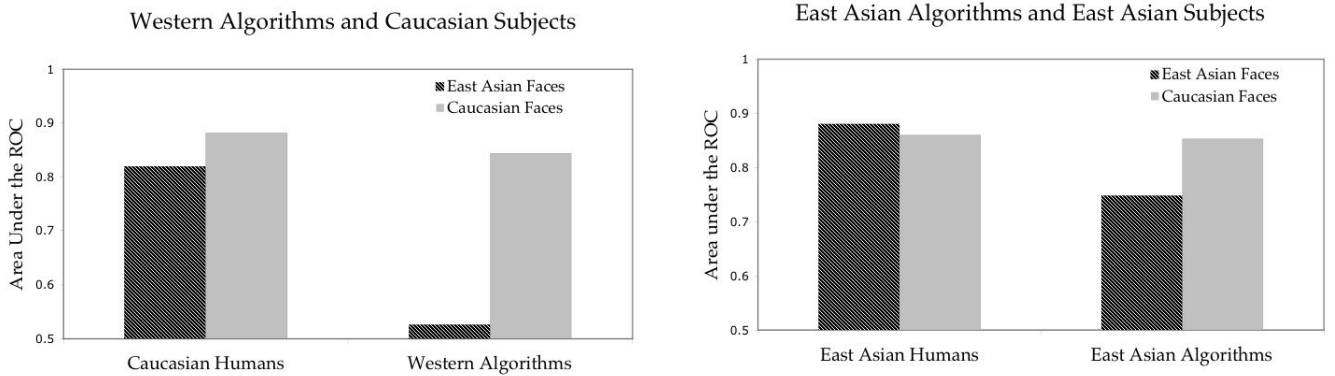


Fig. 3. Performance of humans and algorithms on all 80 pairs of East Asian and Caucasian faces in Experiment 2. Performance is reported with AUC for algorithms and the A' estimate of AUC for humans. The top graph reports performance of Caucasian subjects and the Western fusion algorithm on Caucasian and East Asian faces; the bottom graph reports performance of East Asian subjects and the East Asian fusion algorithm on Caucasian and East Asian faces.

of algorithms robustness. Specifically, algorithm evaluations for suitability in particular application venues should be made using a test population with comparable demographics.

Although it is clear from these experiments that algorithm performance estimates made using populations with different demographics do not converge, understanding the mechanisms behind this finding is challenging. This is due primarily to the fact that the algorithms evaluated in the FRVT 2006 were submitted to the NIST as executables, with no access to source code or to the training sets incorporated by the algorithms during development. Using the human behavioral literature as a guide, and the human findings with this set of test faces, however, we can consider possible causes of the findings that might be common to humans and algorithms.

To understand the complete pattern of results, we start by comparing the other-race effect for the humans and for the two algorithms tested. In all three cases, the effect is defined by an interaction between the race of the face and the race (demographic origin) of the participants (algorithms). A complete crossover interaction was found for the algorithms in Experiment 1 when all available face pairs were tested and when the results focused on low false accept rates. There was also an advantage for East Asian faces. By complete crossover, we mean that the East Asian algorithm was better on East Asian face pairs and the Western algorithm was better on Caucasian faces. This symmetry of crossover defines a “classic other-race effect”. In Experiment 2, humans showed a lopsided interaction, with Caucasian subjects substantially more accurate with Caucasian faces and East Asian subjects slightly more accurate with East Asian faces. The algorithms in Experiment 2 showed an interaction, tilted in favor of Caucasian faces, but with a very large Caucasian face advantage for the Western algorithm and a smaller Caucasian face advantage for the East Asian algorithm. All three results indicate an other-race effect.

In addition to these other-race findings, face race, *per se*, is also an important performance factor in the two experiments. In Experiment 1, performance was reported at low false accept rates using all available face pairs (including some with mis-matched demographics, e.g., gender age). This experiment showed an

advantage for East Asian faces, consistent with a recent study on the FRVT 2006 data where only matched face pairs were considered [22]. Combined, these findings suggest that the low false accept criterion is the primary cause of the East Asian face advantage.

In Experiment 2, where performance was reported over the full range of false accept rates, there is some evidence for a Caucasian face advantage. This was supported both by the Caucasian face advantage seen for the algorithms in Experiment 2 and by the larger other-race deficit Caucasian participants showed in comparison to East Asian participants. This Caucasian advantage could occur potentially if Caucasian faces are inherently more discriminable than East Asian faces. However, data from meta-analyses on human behavioral studies of the other-race effect [4], [3], show no evidence for inherent discriminability differences for faces of different races—making this an unlikely general explanation for the results. For individual test sets, however, there may be small differences in the discriminability of different races of faces. The slight (non-statistically significant) advantage human participants showed for the Caucasian faces combined with the overall advantage of the algorithms on the Caucasian faces in Experiment 2 is consistent with the possibility that the Caucasian faces from this particular data set were inherently easier to discriminate than the East Asian faces.

An equally valid alternative possibility, however, is that both the humans and algorithms had somewhat “more experience” with Caucasian faces than with East Asian faces. As noted, the human participants in this experiment were of East Asian and Caucasian descent, but were recruited from a university in the U.S. in a city where Caucasians comprise the majority of the local population. In fact, most behavioral studies of the other-race effect are conducted with participants of two different races from the same local venue where one of the two races is the local majority race. In these cases, the other-race effect is commonly superimposed on a small local face race advantage.

For algorithms, “experience” refers to the amount and nature of training employed. The relevant component of experience for this study is the extent to which algorithms were trained with different races of faces. On this question, we have no direct knowledge

but we know the following about data availability. All of the research groups included in the fusion algorithms tested here participated in the Face Recognition Grand Challenge (FRGC). The FRGC preceded the FRVT 2006 by two years and one of the data sets used in the FRVT 2006 was collected at the same site as the FRGC data set. The test faces for the FRGC and FRVT 2006 comprised mutually exclusive sets of the faces from the high resolution database developed for these large scale tests [18]. The FRGC data set was composed of a strong majority of Caucasian faces (70%) and a minority of East Asian faces (22%). It is probable that all of the algorithms made use of the FRGC training faces in preparing for the FRVT 2006, and therefore had some experience with Caucasian faces. In addition, training procedures implemented by the East Asian algorithms prior to the FRGC and FRVT 2006 may have included more East Asian faces than training procedures implemented by Western algorithms. Thus, analogous to the East Asian participants in the behavioral experiments, the experience of the East Asian algorithms might have been based on both Caucasian and East Asian faces. Analogous to the Caucasian participants, experience for the Western algorithm training may have strongly favored Caucasian faces.

In conclusion, the performance of state-of-the-art face recognition algorithms varies as a joint function of the demographic origin of the algorithm and the demographic structure of the test population. This result is analogous to findings for human face recognition. The mechanisms underlying the other-race effect for humans are reasonably well understood and are based in early experience with faces of different races. Although our hypotheses about the mechanisms underlying the algorithm effects are still tentative, the effects we report are not. The present results point to an important performance variable combination that has not received much attention. The results also suggest a need to understand how the ethnic composition of a training set impacts the robustness of algorithm performance. Finally, from a practical point of view, recent studies indicate that algorithms are now capable of surpassing human performance matching face images across changes in illumination [28], [25] and on the task of recognition from sketches [29], [30]. This increases the likelihood that face recognition algorithms will find new real-world applications in the near future. In these cases, there is a pressing need to test algorithms intended for applications in venues with highly diverse target populations using face sets that match statistics of the demographics expected in these venues.

ACKNOWLEDGMENT

This work was supported by a contract to A. O'Toole from TSWG. P. J. Phillips was supported in part by funding from the Federal Bureau of Investigation. The identification of any commercial product or trade name does not imply endorsement or recommendation by NIST.

REFERENCES

- [1] R. S. Malpass and J. Kravitz, "Recognition for faces of own and other race faces," *Journal of Personality and Social Psychology*, vol. 13, pp. 330–334, 1969.
- [2] R. K. Bothwell, J. C. Brigham, and R. S. Malpass, "Cross-racial identification," *Personality & Social Psychology Bulletin*, vol. 15, pp. 19–25, 1989.
- [3] C. A. Meissner and J. C. Brigham, "Thirty years of investigating the own-race bias in memory for faces: A meta-analytic review," *Psychology, Public Policy, and Law*, vol. 7, pp. 3–35, 2001.
- [4] P. N. Shapiro and S. D. Penrod, "Meta-analysis of face identification studies," *Psychological Bulletin*, vol. 100, pp. 139–156, 1986.
- [5] D. Levin, "Classifying faces by race: The structure of face categories," *Journal of Experimental Psychology: Learning, Memory and Cognition*, vol. 22, pp. 1364–1382, 1996.
- [6] —, "Race as a visual feature: Using visual search and perceptual discrimination tasks to understand face categories and the cross-race recognition deficit," *Journal of Experimental Psychology: General*, vol. 129, pp. 559–574, 2000.
- [7] A. E. Slone, J. C. Brigham, and C. A. Meissner, "Social and cognitive factors affecting the own-race bias in whites," *Basic & Applied Social Psychology*, vol. 22, pp. 71–84, 2000.
- [8] G. Bryatt and G. Rhodes, "Recognition of own-race and other-race caricatures: Implications for models of face recognition," *Vision Research*, vol. 38, pp. 2455–2468, 1998.
- [9] A. J. O'Toole, K. A. Deffenbacher, D. Valentin, and H. Abdi, "Structural aspects of face recognition and the other-race effect," *Memory & Cognition*, vol. 22, no. 2, pp. 208–224, 1994.
- [10] P. M. Walker and J. W. Tanaka, "An encoding advantage for own-race versus other-race faces," *Perception*, vol. 32, pp. 1117–1125, 2003.
- [11] D. J. Kelly, P. C. Quinn, A. M. Slater, K. Lee, L. Ge, and O. Pascalis, "The other-race effect develops during infancy: Evidence of perceptual narrowing," *Psychological Science*, vol. 18, pp. 1084–1089, 2007.
- [12] C. A. Nelson, "The development and neural bases of face recognition," *Infant and Child Development*, vol. 10, pp. 3–18, 2001.
- [13] S. Sangrigoli, C. Pallier, A. M. Argenti, V. A. G. Ventureyra, and S. de Schonen, "Reversibility of the other-race effect in face recognition during childhood," *Psychological Science*, vol. 16, pp. 440–444, 2005.
- [14] P. K. Kuhl, K. H. Williams, and F. Lacerdo, "Linguistic experience alters phonetic perception in infants by 6 months of age," *Science*, vol. 225, pp. 606–608, 1992.
- [15] N. Furl, P. J. Phillips, and A. J. O'Toole, "Face recognition algorithms and the other-race effect: computational mechanisms for a developmental contact hypothesis," *Cognitive Science*, vol. 26, pp. 797–815, 2002.
- [16] H. Moon and P. J. Phillips, "Computational and performance aspects of PCA-based face-recognition algorithms," *Perception*, vol. 30, pp. 303–321, 2001.
- [17] R. Gross, S. Baker, I. Matthews, and T. Kanade, "Face recognition across pose and illumination," in *Handbook of Face Recognition*, S. Z. Li and A. K. Jain, Eds. Springer, 2005, pp. 193–216.
- [18] P. J. Phillips, P. J. Flynn, T. Scruggs, K. W. Bowyer, J. Chang, K. Hoffman, J. Marques, J. Min, and W. Worek, "Overview of the face recognition grand challenge," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2005, pp. 947–954.
- [19] P. J. Phillips, H. Moon, S. Rizvi, and P. Rauss, "The FERET evaluation methodology for face-recognition algorithms," *IEEE Trans. PAMI*, vol. 22, pp. 1090–1104, October 2000.
- [20] P. J. Phillips, P. J. Grother, R. J. Micheals, D. M. Blackburn, E. Tabassi, and J. M. Bone, "Face recognition vendor test 2002: Evaluation report," National Institute of Standards and Technology, Tech. Rep. NISTIR 6965, 2003, <http://www.frvt.org>.
- [21] G. H. Givens, J. R. Beveridge, B. A. Draper, P. J. Grother, and P. J. Phillips, "How features of the human face affect recognition: a statistical comparison of three face recognition algorithms," in *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2, 2004, pp. 381–388.
- [22] J. R. Beveridge, G. H. Givens, P. J. Phillips, B. A. Draper, and Y. M. Lui, "Focus on quality, predicting FRVT 2006 performance," in *Eighth International Conference on Automatic Face and Gesture Recognition*, 2008.
- [23] P. J. Phillips, W. T. Scruggs, A. J. O'Toole, P. J. Flynn, K. W. Bowyer, C. L. Schott, and M. Sharpe, "FRVT 2006 and ICE 2006 large-scale results," National Institute of Standards and Technology, Tech. Rep. NISTIR 7408, 2007. [Online]. Available: <http://iris.nist.gov/>
- [24] P. Grother, "Face recognition vendor test 2002: Supplemental report," National Institute of Standards and Technology, Tech. Rep. NISTIR 7083, 2004, <http://www.frvt.org>.
- [25] A. J. O'Toole, P. J. Phillips, and A. Narvekar, "Humans versus algorithms: Comparisons from the FRVT 2006," in *Eighth International Conference on Automatic Face and Gesture Recognition*, 2008.
- [26] N. A. Macmillan and C. D. Creelman, *Detection theory: A user's guide*. Cambridge: Cambridge University Press, 1991.
- [27] I. Pollack and D. Norman, "A non-parametric analysis of recognition experiments," *Psychonomic Science*, vol. 1, pp. 125–126, 1964.
- [28] A. J. O'Toole, P. J. Phillips, F. Jiang, J. Ayyad, N. Pénard, and H. Abdi, "Face recognition algorithms surpass humans matching faces across

- changes in illumination,” *IEEE Trans. PAMI*, vol. 29 1642-1646, pp. 1642–1646, 2007.
- [29] X. Tang and X. Wang, “Face sketch synthesis and recognition,” in *Proc. Ninth IEEE Int’l Conf. Computer Vision*, 2003, pp. 687–694.
- [30] —, “Face sketch recognition,” *IEEE Trans. Circuits and Systems for Video Technology*, pp. 50–57, 2004.